

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Application of:

Johan Ageby et al.

Application. No.: 10/025,991

Filed: December 26, 2001

Title: METHODS FOR CONTROLLING RESOURCES IN A COMMUNICATION
NETWORK



Art Unit: 2661

RECEIVED

MAR 22 2002

Technology Center 2600

CLAIM FOR PRIORITY

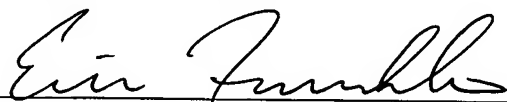
Assistant Commissioner for Patents
Washington, DC 20231

Sir:

Certified copies of corresponding Swedish Application No. 0004896-7, filed December 29, 2000, and Swedish Application No. 0102471-0, filed July 10, 2001, are attached. It is requested that the right of priority provided by 35 U.S.C. 119 be extended by the Patent Office.

Respectfully submitted,

Date: March 20, 2002


Eric J. Franklin, Reg. No. 37,134
Swidler Berlin Shereff Friedman, LLP
3000 K Street, NW, Suite 300
Washington, DC 20007
Telephone: (202) 424-7500

PRV

PATENT- OCH REGISTRERINGSVERKET
Patentavdelningen



Intyg Certificate

RECEIVED

MAR 22 2002

Technology Center 2600



Härmed intygas att bifogade kopior överensstämmer med de handlingar som ursprungligen ingivits till Patent- och registreringsverket i nedannämnda ansökan.

This is to certify that the annexed is a true copy of the documents as originally filed with the Patent- and Registration Office in connection with the following patent application.

(71) Sökande Net Insight AB, Stockholm SE
Applicant (s)

(21) Patentansökningsnummer 0102471-0
Patent application number

(86) Ingivningsdatum 2001-07-10
Date of filing

**CERTIFIED COPY OF
PRIORITY DOCUMENT**

Stockholm, 2002-01-02

För Patent- och registreringsverket
For the Patent- and Registration Office

Christina Vängborg
Christina Vängborg

Avgift
Fee 170:-

jw/

Ref.: 55324 SE

Applicant: Net Insight AB**Dynamic Protocol****5 THE BACKGROUND OF THE INVENTION AND PRIOR ART**

The present invention relates generally to a telecommunications protocol, which provides a mechanism for negotiating resources between interfaces in the system. More particularly the invention relates to a method of allocating resources in a synchronous
10 time division multiplex communications system according to the preamble of claim 1 and a communications system according to the preamble of claim 11. The invention also relates to a computer program according to the preamble of claim 9 and a computer readable medium according to claim 10.

15 The known protocols for allocating resources in communications system where the transmission resources are constituted by time slots in a repeating frame structure have presumed that the resources are allocated at initial configuration of the system. Such procedure requires substantial planning by the network
20 operator and is very inflexible to later alterations of the system's topology and/or changes in the resource requirement from the various interfaces in the system.

Furthermore, there is a risk that the system is configured such that an overlap may occurs, between the configured resources.
25 For instance, the number of configured resources may not correspond to the number of actually available resources.

SUMMARY OF THE INVENTION

The object of the present invention is therefore to provide a resource allocating solution, which alleviates the problems above and thus offers a simple and adaptable distribution of resources in a system of any size.

- 5
- According to another aspect of the invention the object is achieved by a method of allocating resources in a synchronous time division multiplex communications system, as initially described, which is characterized by the following steps;
- 10 sending a link status message from an interface whenever the interface registers a change in the topology of the system, sending a gather message from an interface whenever the interface requests a revision of a current ownership distribution of resources, sending a sync message from the master interface
- 15 as an indication of a current distribution of ownership with respect to the resources between the interfaces in the system, and, for each interface, generating a distribution of the ownership to the resources on basis of the interface's topological position and a latest received sync message.
- 20 According to a further aspect of the invention the object is achieved by a computer program directly loadable into the internal memory of a computer, comprising software for performing the above proposed method when said program is run on a computer.
- 25 According to another aspect of the invention the object is achieved by a computer readable medium, having a program recorded thereon, where the program is to make a computer perform the proposed method.
- 30 According to one aspect of the invention the object is achieved by a communications system as initially described, which is characterized in that it comprises at least one node, which in turn includes one or more of the interfaces. The node is presumed to be adapted to effect the proposed method.

- The invention offers an efficient, reliable and fair solution for dynamically allocating transmission resources in a communications system. The proposed solution on one hand makes manual configuration unnecessary. On the other hand, if a
- 5 system still is manually configured, the invention safeguards against any erroneous or conflicting configurations. This is, of course a very desirable feature from a network operator's point-of-view.

BRIEF DESCRIPTION OF THE DRAWINGS

- 10 The present invention is now to be explained more closely by means of preferred embodiments, which are disclosed as examples, and with reference to the attached drawings.

Figures 1a-d show different examples of allocation domains to which the invention is applicable,

- 15 Figure 2 illustrates a so-called short-circuit scenario,
- Figure 3 illustrates a typical probe according to an embodiment of the invention,
- Figure 4 illustrates a typical borrow and return of resources according to an embodiment of the invention,
- 20 Figure 5 shows important transitions of a quark machine according to an embodiment of the invention,
- Figure 6 shows an example of messages sent when an interface is appended to a bus according to an embodiment of the invention,
- 25 Figure 7 shows a first example of messages sent when a ring is closed according to an embodiment of the invention,
- Figure 8 shows a second example of messages sent when a ring is closed according to an embodiment of

the invention, and

Figures 9a, b show two important cases in which resources are allocated according to embodiments of the invention.

5 DESCRIPTION OF PREFERRED EMBODIMENTS OF THE INVENTION

This proposed solution relates to a mechanisms for negotiating resources between interfaces in a synchronous time division multiplex communications system having a master interface
10 communicating with one or more slave interfaces, and in which the resources between the interfaces are represented by time slots in a repeating frame structure. Thus, the solution may be applied in a system of dynamic synchronous transfer mode (DTM) type. In such system, the solution can be accomplished
15 by the Resource Management Protocol (DRMP). DRMP is a token passing mechanism for negotiating between interfaces which resource units that are available. It is divided into two orthogonal mechanisms, ownership and access.

Definitions

20 **allocation domain:** the same as a bypass chain but if the topology is point-to-point or bus, the last node is not counted as member of the AD.

access right: the right to write on a certain slot on an interface for a certain number of bypass hops

25 **access token:** a token that is passed around to grant access right between interfaces on an allocation domain

bool: a logical variable type, which can take two values, true or false

30 **fairness algorithm:** the algorithm used to determine actual ownership ranges from a set of requested policies from the interfaces in the AD

master interface: the interface in the AD having the lowest mac address.

ownership: the obligation to manage an access token with respect to lending and issuing probe and kill messages

quark: the smallest resource unit it is one slot wide and one bypasshop long

- 5 **topology:** a set of two or more interfaces connected in a bypass chain that is either closed or open.

Abbreviations

For the purpose of the present document, the following abbreviations apply:

10	AD	Allocation Domain
	BR	BitRate
	DCC	DTM Control Channel
	Distown	Distribute Ownership
	DLC	Data Link Change
15	DLSP	DTM Link State Protocol
	DRMP	DTM Resource Management Protocol
	DO	Dynamic Ownership
	Mac	MAC Address.
	Msg	Message
20	Prrpy	probe reply
	Qreq	quark request
	Qret	quark return
	Qt	quark transfer
	Ptp	Point-to-point

25 Problem domain

DRMP efficiently distributes write access rights to the time slots in a DTM system. This is done with the following mechanisms:

- Ownership and Announcement of resources.
- Borrowing and Probing of resource tokens.

30 Resource ownership

Consider the everyday concept of ownership. Ownership does not automatically grant access right. Remember that even if something is owned by a certain unit, some other unit might have borrowed it and thus made it inaccessible to the owner.

Static ownership

In this mode, all the resource units have had their ownerships defined at startup time. It is still possible to change the ownership distribution of a running system. This requires manual
5 intervention by the operator, often in more than one unit, something that is very complicated.

Dynamic ownership

In this mode, ownership is negotiated by letting all units initially state how many resources they want. The master interface then
10 distributes these policies to all the participating interfaces, including the master itself. Dynamic mode eliminates the necessity for configuration by the operator. Overlapping may occur transiently but the DO ensures that this state does not persist and that it is never dangerous.

Resource announcement

15 Resource announcement is done by each unit telling all the other units how many tokens it is willing to lend out. When a unit receives an announcement it stores that information. This is used when borrowing.

20 Later it uses this information to decide which units to borrow from.

Borrowing

If a unit finds that it has not got enough resources it will attempt to borrow access rights to the other units resources. This is
25 done by passing tokens between the borrower and the most suitable owners according to the announce tables. When the tokens are to be returned they are always returned to the owner.

Probe / Kill

The Probe mechanism has two main objectives:
30 - Recreate lost resource tokens or initially create them.
- Resolve situations when two units simultaneously claim to have the access right to a resource.

The probing is carried out in the same way for both cases. In order to do a successful probe of a set of tokens, the owner of
35 the resource must make one query to all other units sharing this

- token and get a reply from all of them. There are three main outcomes of a successful probe. The token may be non-allocated, allocated once or doubly allocated. Probing involves sending many messages across the network, especially at bootstrap. This can take considerable time. A useful special case is the point-to-point topology, in which the two units create the resources locally.

Fault handling

- The fault situations and what DRMP does to fix them when they have occurred are described below.

Borrowing failed

- Since the number of resources announced by a unit only is valid at the time of the announcement, borrowing does not always succeed. When a unit has attempted to borrow tokens and have not managed to find the minimum required.

Ownership overlapping

Ownership overlapping causes a problem in the static ownership mode only. The dynamic ownership protocol ensures that either the ownerships are non-overlapping or the probe is not active.

Loss of access tokens

- The most common case of lost access tokens would be at startup since units initially are without resources and then acquires them by probing. Access tokens can also be lost by message losses on the control channel because of congestion, buffer overflows etc.

Doubly allocated tokens

- This is the most serious case of faults in DRMP since it implies that two or more units have write access to the same resources at the same time. In this case, integrity of the data transported in the channels using these resources is violated. DRMP has been designed with the main objective of never allowing this to happen. It has been shown however that if a number of factors work together, it might nevertheless occur. Therefore, the kill mechanism exists to allow recovery from those unusual

situations. An important note to make here is that even if the owner has the resource in use locally it must probe periodically anyway, since someone else may wrongly be using that same resource.

5 **Allocation domain**

10 An allocation domain is a set of interfaces connected in sequence with each other. The allocation domain type can be either of point-to-point, bus or ring. In the bus cases, the last interface is defined not to be part of the allocation domain. Its transmitter is not used and hence it needs no outgoing resource. The allocation domain interpretation of a DLSP topology is also what is received via a DLC message. See figures 1a - d for examples. The figures 1a - d show: a five-node bus, ring, two-node bus and point-to-point topologies respectively.

15 **Short circuit**

20 Two or several interfaces may belong to the same node. This is sometimes referred to as "short-circuited" interfaces. Figure 2 shows an example of this. This is handled the same way as if C1 and C2 would have been on different nodes, with the exception that C1 and C2 communicate internally in the node and that channels originating in C1 may not pass beyond C2 and v.v.

25 A short-circuit is two interfaces that both belong to the same node and to the same allocation domain. In this case the ports C1 and C2 "speak" to each other using the ordinary DRMP messages, via a node-local control channel.

30 **Probe**

The probe is the mechanism responsible for detecting that there is a resource token missing and if so, recreate that token locally. Probing only takes place for slots that we consider owned. The probe always asks all the other nodes in the allocation domain if they use the resource. If not, then the resource token is recreated. The probe is the only mechanism that will recreate tokens (other mechanisms are only responsible for transferring tokens around). To avoid the risk of ownership overlaps and

thus the risk of double bookings, the probe is turned off during ownership changes.

Figure 3 shows a typical probe session. It is a probe that just checks that a borrower still has the slot. Probing is done

5 because of two reasons:

- At bootup to get the resources initially, remember that all nodes must ask around for its resources since they may be borrowed out from a previous owner.
- Tokens can get lost when transmitted from one node to another.
- 10 - The node has just started. It always starts with no resources but with an ownership range assigned. This means it also has a responsibility to probe its resource tokens at a regular interval. At bootup, the probe is intensified, (sometimes called the turbo probe). The turbo strategy is to make at least one successful probing for each resource unit and when that is done go to
- 15 probing at lower rate (and thus spend less CPU). The turbo is also re triggered whenever the topology of the specific allocation domain changes (due to fiber or node failures).
- Messages get lost on networks. The reasons are many, congestion in buffers or bit errors on the transmission media. The assumption is however that these occasions are rare and thus
- 20 we settle for a relatively slow mechanism for the detection of lost resources.

Borrowing and lending of token(s)

25 Tokens can be borrowed from other nodes in the allocation domain. In figure 4 we see a typical borrow session. Interface A1 borrows from B1. After some time the channel is torn down and the borrowed resource is no longer needed and is thus returned.

30 Ownership of token(s)

The ownership of tokens is a process orthogonal to the access right. Changing ownership does not necessarily change the access right and vice versa. The intention of the distributed ownership defined in this protocol is to try to have the resources

35 where they are needed, since borrowing and lending yields a

higher cost resource-wise than a simple local allocation at the local interface that needed the resource.

Access to token(s)

- 5 The access to tokens means that we either are already using the resource or have the right to do so, for instance by holding it in the local free list or using it for an active channel.

Gather

- 10 The message sent out from an interface that is "unhappy" with the current ownership distribution, or rather, the last received Sync message, (see Figure 6 and forth for examples), did not give the same limit for this interface as what we have locally.

Sync

- 15 A message initiated from the master node only or a node which thinks it is master, aiming at having all interfaces have the same idea of what resources they own.

Bootup

The process taking place when the power is switched on.

Master interface

- 20 The master interface is defined as the node having the lowest mac address. At any given time there is no guarantee that the nodes have consistent topology info.

Transitional master interface

- 25 This is the name of an interface, which becomes the master for a short period, since it has not yet got the correct topology information. Other interfaces may have other info on the topology and the mastership. Since the full allocation domain information is contained in the sync message, it is quite easy for an interface to determine that a sync message should be ignored.

- 30 Probe

A mechanism aiming at doing consistency checks for the allocation domain. It is also responsible for resolving resource conflicts and re-creating tokens that have been lost or, at boot time, do an initial create of the resources.

Topology

A Topology is a set of interfaces given to us as an ordered list from DLSP. A topology is either ring or bus and the last terminating interface is included in the topology.

5 Fairness algorithm

This is the algorithm that is used by the master of the dynown system to calculate the total amount of slots that each node should have. Please note that it is the policing parameters that are sent from each node which are distributed to each node and
10 the calculation of the fairness takes place at each node, not only at the master. Therefore it is important that the fairness algorithm is defined in a precise manner.

Range change

This is an incoming event to the executive dynown telling it how
15 many slots a specific interface would like to have. If the requests for the total amount of slots exceed the total available range, a fairness algorithm is used for this calculation.

Double booking

Double booking is said to have occurred when two or more
20 interfaces believe they have access right to a certain resource unit or set thereof. There are two cases of this, disastrous and potentially disastrous. The disaster is defined as at least one slot being used for a channel on at least two nodes when the scope of the slot(s) are overlapping.

25 Intrinsic message

An intrinsic message, as opposed to an incoming or outgoing message, is a message that goes from one instance of a module on one interface to another instance of the same module in an other interface. The probe is a good example of an intrinsic
30 message, since it is only the respective DRMP instances in two or more nodes that talk to each other. See figure 3 and 4 for examples of DRMP-intrinsic message passing.

Quark

This is the smallest resource unit available. It is defined to be
35 one slot "high" and one physical link "wide". The quark machine

defines resource management with the assertion that any implementation will be an aggregation of several of these quarks in bulk for various optimizations. When uncertain how to handle a special case or event in the implementation, refer to the
5 behavior of the quark machine.

Mine

An event telling the quark machine it now owns its quark.

Not mine

10 An event telling the quark machine it does no longer own its quark.

Alloc

Telling the quark machine its resource is used for a channel.

Dealloc

15 Telling the quark machine its resource is no longer used for a channel.

Fragmentation

20 This is pretty much the same type of fragmentation that occur in computer disks and memories. Several small channels are first allocated consecutively. Deallocation is not done consecutively; hence we end up with fragments of slots that require more computation and memory.

Worst case fragmentation

25 Sometimes the fragmentation of the resources has to be limited. The reasons may be many. The worst case of fragmentation is the case where the busy/free resource map resembles a checkerboard. It is used to calculate maximum message sizes though no formal proof exists that it generates the largest possible messages.

Bitrate

30 This is a measure of the capacity of a channel or a link. It is measured in slots. A slot is defined to be 512 kb/s. For example: 200 slots are about 100Mb/s of bitrate.

DETAILED PROTOCOL DESCRIPTION

The Quark machine

The Quark machine is a model in which we assume we only have one resource unit moving around in the system. The idea is that this simplifies initial modeling and verification. It is also assumed that the generalization to larger resource units is straightforward since they can be handled in blocks.

Simplified graphical diagram

Here is the graphical description. Please note that only the "normal" transitions are shown in this diagram. For all events on all states, see the below tables.

Events

- QReq (Quark Request) - This message is a request for a token it does not in it self change the state of the global system.
- QT (Quark Transfer) - This means that a token is being sent from one interface in the allocation domain to another. This token is also said to belong to a session. That is, it is aimed at a specific channel in the borrowing node. There are two possibilities: The resource message reaches the destination and the two state-machines changes state or the resource message gets lost for some reason. Only the state machine where the message leaves changes places.
- Qret (Quark Return) - Similar to QT, but this is just returned to another interface in the allocation domain and left in the free resource pool at the other node.
- Mine - This represents a change of ownership such that the local interface that gets the message now owns the resource.
- NotMine - This represents a change of ownership such that the interface that gets the message is no longer responsible for managing the resource.
- Alloc - This is a local request that should only occur on a Free resource. Most implementations will probably not be able to do anything else, since the resources are taken from a pool of available resources. The resource is allocated and moves to state Busy.
- Dealloc - A local resource is returned to the pool of Free resources.

- Probe - This message is used to ask other interfaces if they are using a certain resource.
- PrRpy - This message is the reply for the Probe, which tells us the state of the resource at that specific interface.
- 5 - Timeout - Currently, this only happens in the Probe state. (Remember that sending a QReq does not change the state of the Quark machine at the node that initiates the message).

States

- Free - Resource is free to use.
- 10 - Lent - Resource is used by someone else.
- Gone - We don't care where the resource is since we don't own it and don't use it.
- Borrowed - We don't own the resource, but we are using it.
- Busy - The resource is ours and we are using it.
- 15 - Probing - The resource is undergoing an examination of whether someone is using it or not.

State transition tables

- This is all the states and events possible for the quark machine. The side effect "WARN" tells us that something has happened that "shouldn't". This means either one of two things:
- 20 - The event has become "impossible" due to implementation choices.
 - The event is "illegal" in that it changes the system to a possibly dangerous, inconsistent or unknown state.

25 Free

The free state represents the fact that the resource is known to be available for use (either remote or local), that is, we *know* that no one else in the system is using the resource and that we are not using it ourselves.

30

Event	Condition	NextState	Action
QReq		Lent	send QT
QT		Free	WARN
QRet		Free	WARN
Mine		Free	
Not Mine		Gone	
Alloc		Busy	
Dealloc		Free	WARN
Probe		Free	
PrRpy		Free	
Kill		Lent	
time-out		--	

Table 1: Transition table for the Free state.

Busy

This state represents the fact that we know we are using the resource locally.

5

Event	Condition	NextState	Action
QReq		Busy	
QT		Busy	WARN
QRet		Busy	WARN
Mine		Busy	
NotMine		Borrowed	
Alloc		Busy	
Dealloc		Free	WARN
Probe		Busy	
PrRpy		Busy	
Kill		Busy	send Dcp Remove
time-out		--	

Table 2: Transition table for the Busy state.

Lent

This state represents the belief that the resource is in use somewhere else in the system and not at our node. The reason this is a belief is because this state represents the state of the system the last time we looked. Remember that in a distributed system things may have changed recently in other units without us knowing. This state is also the start state for resources

10

owned by the local node. At start we assume that someone is using the resource until we have asked everyone concerned.

Event	Condition	NextState	Action
QReq		Lent	
QT		Lent	WARN
QRet		Free	
Mine		Lent	
NotMine		Gone	
Alloc		Lent	WARN
Dealloc		Lent	WARN
Probe		Lent	
PrRpy		Lent	
Kill		Lent	
time-out		Probing	Send Probe's

Table 3: Transition table for the Lent state.

5 Gone

This represents the fact that we do not consider the resource to be ours and that we do not have it borrowed right now.

Event	Condition	NextState	Action
QReq		Gone	
QT		Borrowed	
QRet		Gone	
Mine		Lent	
NotMine		Gone	
Alloc		Gone	WARN
Dealloc		Gone	WARN
Probe		Gone	send PrRpy (Gone)
PrRpy		Gone	
Kill		Gone	
time-out		—	

Table 4: Transition table for the Gone state.

Borrowed

This means we have borrowed in a resource that someone else owns and that it is currently in use at our node.

Event	Condition	NextState	Action
QReq		Borrowed	
QT		Borrowed	WARN
QRet		Borrowed	WARN
Mine		Busy	
NotMine		Borrowed	
Alloc		Borrowed	WARN
Dealloc		Gone	send QRet
Probe		Borrowed	send PrRpy (Borrowed)
PrRpy		Borrowed	
Kill		Borrowed	send Dcp Remove
time-out		—	

5 Table 5: Transition table for the Borrowed state

Probing

This means that we have decided to ask for the resource but we have not yet obtained answers from all or any of the other nodes involved.

Event	Condition	NextState	Action
QReq		Probing	
QT		Probing	WARN
QRet		Free	
Mine		Probing	
NotMine		Gone	
Alloc		Probing	WARN
Dealloc		Probing	WARN
Probe		Probing	
PrRpy	Not Last PrRpy	Probing	p[i]=state
Kill		Probing	
time-out		Lent	
PrRpy	Last && InUse(p, state)	Lent	p[i]=state
PrRpy	Last && UnUsed(p, state)	Free	p[i]=state

Table 6: Transition table for the Probing state.

The functions inUse and unUsed deserve a special explanation. When we apply these functions they depend on the sub condition that the last probe reply has just arrived and that p contains all but the last reply. Remember that the condition is checked *before* we take action.

- InUse - This function is defined as true if *any* of its arguments are true. If more than one of its arguments are true a protocol error has occurred. This should trigger the kill function. The unUsed function implies that *all* of p is false (i.e. unused).
- UnUsed - This function is true if none of its arguments are true. This tells us that a recreation of the resource should be done.

The Dynamic ownership machine

The DO machine is the system for ownership negotiations. The states Free, Lent, Busy and Probing mean we have ownership rights to a resource. The state Borrowed or Gone means someone else has. The aim of this protocol is to negotiate ownerships in such a way that they never overlap before turning on the probe.

Typical scenarios

- 5 The aim of the design of dynamic ownership distribution is to cater for all situations, common or uncommon. An attempt has been made to optimize for some main scenarios, such as minimizing the time for the bootup sequence. See the figures 6, 7 and . Either the master gets the DLC (from DLSP) last, or some other node (a slave) will not have gotten it. Remember that in a distributed system, it is impossible to tell what goes on in the other nodes.

10 Transition diagrams

This is the transition diagram for the executive part of the dynamic ownership distribution.

State types

- 15 The following types are used for state variables in the DO machine:

Type	Description
AD	A representation of an ordered set of mac addresses
Bool	see Abbreviations
Integer	An unsigned number at least 32 bits large.
SyncStates	A variable, which can take two values: Idle or Wait.
Interface	This is just a 48 bit mac
List<T>	Some kind of ordered template which can hold several of these variables.
Limit	The Limit is a structure with two attributes: An ownership limit and a starvation flag, Integer and Bool respectively.
StarvFlag	A flag that indicates if an interface is set to starvation mode or not (i.e. it does not accept any new channels to/from or over it).

Table 7: The types used in the DO machine

State variables

The following state variables exist for each local port on the local node:

Name	Type	Description
pending	AD	Holds the topology info pending before synchronization is done. Initial value should be the own interface if a range_change triggers or the DLC topology if DLC triggers.
sync State	Sync States	This holds the state we are currently in, Idle or Wait. Idle means we are happy and ready to issue and reply to probes, however we claim that all DRMP links should be Idle before the global isMaster call yields true. Wait means that we are waiting for a synchronization and thus does not probe on any of our interfaces. This is a security measure since an interface might suddenly become part of another link when we start to move fibers around.
me	Interface	This is the Interface identifier of the local interface, we should have this as one of the members in the active and pending lists. It is to be regarded as a constant for a given state machine.
limits	List <Limits>	This holds the pending ownership ranges. We compare this to the incoming Gather's and determine if a new sync should be done. This is done if the limits change. If we are not the master we only store the limits. The limits are to be assigned a default zero at first, stating that that specific node does not need resources. When that node receives the sync, it should send a Gather, with its real number, protesting in order to get resources.
mylim	Limit	This holds the ownership limit. The own limit that is set only by RangeChange message. Initially one element of zero for the local Unit, or, if this State machine instance is booted with range_change an initial element with the given limit.
stbit	StarvFlag	This is a flag, which is given to us via operator interface. Then it is sent with the gather message, the same way as with the limit. All implementations need not be able to do this, but all must handle the drop of resources when a starvation bit is received via the sync message.
stbits	List <Starv Flag>	This holds flags for all interfaces in the allocation domain, that have requested that any free resources crossing their physical link must be dropped, this is used in the sync message.

Table 8: The state variables on a per local interface basis.

Queries

Queries are typically implemented as synchronous function calls in local software.

Probe conditions

- 5 This message is a query from executive DRMP concerning whether or not it is allowed to probe or reply to probes right now. The criteria is as follows:
- For each known local interface, the state is currently set to idle. This means it is allowed to probe right now.
- 10 - One, several or all of the local interface states are set to Wait. This means it is not allowed to probe right now for *any* of the interfaces.

- Implementations might gain from caching the value of this variable and only recalculate it on any event that alters the state of any interface and changes to/from Wait/Idle.
- 15

Incoming Message

Asynchronous events going into the DO machine.

DLCDescription

- 20 DLC(Data Link Change) comes originally from DLSP. When stored, it is to be regarded with the probeUnit view (one less interface if buss or ptp).



master(r)	state	Action	New State	Send
no	Idle	pending:=r, clear(oldlimits), limits[me]:=myl im	Wait	Gather(me, mylim)
no	Wait	pending:=r, clear(oldlimits), limits[me]:=myl im	Wait	
yes	Idle	pending:=r, clear(oldlimits), limits[me]:=myl im	Wait	Sync(me, r, limits)
yes	Wait	pending:=r, clear(oldlimits), limits[me]:=myl im	Wait	Sync(me, r, limits)

Table 9: Reception of DLC(Interface me, AD r)

- Clear(oldlimits) means that we remove the entries in our limits-container for possible interfaces that have been removed from the topology, (as given by the DLC event), but we want to keep interfaces already there.

RangeChange

Description

- This is done from the user interface of a specific node (resedit & friends). In the case of the master getting a range change we send the gather internally to "ourselves".

Action	Send
mylim:=l	Gather(me, mylim)

Table 10: Reception of RangeChange(Limit l)

- If this message gets lost on the way to the master, we rely on periodic retransmission of Gather. If this interface is the master, the Gather is sent locally.

Intrinsic Messages

Intrinsic messages are message which talks from one instance of DO machine in one task to another in another task, (usually another code-executing unit).

5 Sync

Description

10 This is the sync message. It is threaded between the interfaces in order from Master to most downstream, then to most upstream and downstream back to the master. Different interfaces may have other info on the topology and the mastership. Since the full allocation domain information is contained in the sync message, it is quite easy for an interface to determine that a sync message should be ignored.

master (pending)	pending = r	mylim= l[src]	State	Action	New State	Send
X	no	X	Idle		Idle	
X	no	X	Wait		Wait	
no	yes	yes	Idle	limits:=l	Idle	DistOwn(l)Sync (me, r, l)
no	yes	yes	Wait	limits:=l	Idle	DistOwn(l)Sync (me, r, l)
no	yes	no	Idle		Idle	Gather(mylim)
no	yes	no	Wait		Wait	Gather(mylim)
yes	yes	yes	Idle	limits:=l	Idle	
yes	yes	yes	Wait	limits:=l	Idle	DistOwn(l)
yes	yes	no	Idle		Idle	Gather(mylim)
yes	yes	no	Wait		Wait	

15 Table 11: Reception of Sync(AD r, List<Limit> l)

Please note that the limits-parameter in the Sync message should be the set of limits given from the interfaces. The actual ownership distribution should be calculated at each node, this of course implies that this function must be the same everywhere.

20 For each interface calculate the start and range, then pass it on to executive DRMP via the RangeChange event. This is because we want each interface to compare the Synced limits to its local mylim. If not happy, the slave sends a Gather with the proper

- value. The reason for doing so is that several desired limit-vectors maps to a certain set of ownership distributions. For example, if there are 5 interfaces and 10 slots total, then setting the policing parameters of all nodes to the same value, will yield the same results (2,2,2,2,2) for the ownership limits, thus we couldn't know whether to send a Gather without our "new" limit or not.

- The ring and bus figures 9a and 9b show the signal paths from the master interface (marked with a crown) C1. In the bus case the sync-messages goes C1-D1-A1-B1-C1, while in the ring case, we go C1-D1-A1-B1-C1. Please also note that even though interface e1 is present it is *not* part of the ownership negotiation.

Gather

- 15 Description

This is periodically sent to the master of a pending link. A 'yes' in the second column, implies that the range for this unit has already been transmitted.

master (pending)	limits[src] = I	State	Action	New State	Send
no	X	Wait		Wait	
no	X	Idle		Idle	
yes	no	Idle	limits[src]:=I	Wait	Sync(me, pending, limits)
yes	no	Wait	limits[src]:=I	Wait	Sync(me, pending, limits)
yes	yes	Wait		Wait	
yes	yes	Idle		Idle	

- 20 Table 12: Reception of Gather(Interface src, Limit I)

Outgoing Messages

DistOwn(List<Limit> I)

DistOwn can be issued several times, with the own-limits (I) changing for each I. The Probe must be off when sending this. If

any StarvFlag bits are set we must drop any free resources that we have accessright to on the scope of that interfaces physical link.

Time-outs

5 (Re) send Gather

The Gather is sent periodically so that even if a message is lost we will eventually get through to the master. This time-out should be in the order of several seconds, since it is transmitted on a regular basis.

10

master(pending)	Action
no	send Gather(me, mylim)
yes	

Table 13: Timeouts

Message formats

This describes the various messages in DRMP. Note that all fields are to be network order. Please also note that the actual
15 formatting on the outgoing stream should be from right-to-left and top-to-bottom.

In some messages, we use the division sign "/". This is to be interpreted as an integer division, i.e. any fraction is immediately truncated by the operation itself. We also use
20 modulo "%" which is interpreted as the remainder of an integer division.

Bits and pieces

This is repetitive parts of messages that are frequent in the actual messages below. They are not themselves to be regarded
25 as complete messages.

Mac addresses field

The Mac address is always placed "to the right" in a DTM frame, with the layout below:

30

0	Mac address
---	-------------

Fields	Slot #	BitVec	Size	Description
Mac address	any	[47:0]	48	This is the field for the 48 bit mac address

Table 14: Mac address fields

Slot fragment

This pattern is common throughout the document, it is the way to transport tokens through the network.

5

amt n	start n	amt n+1	start n+1
-------	---------	---------	-----------

Fields	Slot #	BitVec	Size	Description
amt n	any	[63:56]	8	Amount of these slots
start n	any	[55:32]	24	Start of this slot fragment
amt n+1	any	[31:24]	8	Amount of next slots
start n+1	any	[23:0]	24	Start of next slot fragment

Table 15: Slot fragment fields

Long slot fragment

- 10 If both amt fields have the value zero slot fragments have a slightly different meaning, this is shown below:

0	start	0	amount
---	-------	---	--------

15

Fields	Slot #	BitVec	Size	Description
start	any	[55:32]	24	Start of this slot fragment
amount	any	[23:0]	24	Amount of slots

Table 16: Long slot fragment fields

This slot fragment representation is more optimal for very large consecutive slot tokens.

Generic header

- 20 This is the generic header for DRMP. Please note that the destination below is the interface we are talking *about*. This is not necessarily the same as the interface used to receive the actual message on the network (Checkout a double bus for an

example of this, it becomes even more obvious when we have more than one interface in an Allocation domain on the same node).

5	Vr	Cmd	0	Destination mac Address
---	----	-----	---	-------------------------

Fields	Slot #	BitVec	Size	Description
Vr	0	[63:61]	3	The version of the protocol (currently 0)
Cmd	0	[60:56]	5	The command of the protocol
Destination mac Addresses	0	[47:0]	48	Destination mac that we are talking about

Table 17: Generic header fields

10 **Statistical information**

These are messages that do not alter any state in the quark machine. They are merely used as information. In a distributed environment they only hold the right data if the system is in steady state. If we do employ dynamic ownership the message with code 2 is used. If we do not have dynamic ownership distribution, we use the code 3 and stuff the ownership start and range in the message as well. All implementations should handle both message types and if an announce with code 2 is received when employing static ownership, the fields Own Start and Own Range should be interpreted as being zero although they are not present in the actual message.

Resource Announce (Dynamic ownership)

This is used if the interfaces ownership distribution is dynamic

0	2	N	broadcast mac addr
0			Announce Source mac
0		br hw 1	Up stream port link layer address 1
br lw 1			Down stream port link layer address 1
0		br hw n	Up stream port link layer address n
br lw n			Down stream port link layer address n
0		br hw N	Up stream port link layer address N
br lw N			Down stream port link layer address N

5

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56][47:0]	56	Header, Cmd=2
N	0	[55:48]	8	This is the number of fragments in this announce message
Announce Source mac	1	[47:0]	48	This is the mac of the announcing interface
br n	2n, 2n+1	2n[55:48] (2n+1)[63:49]	24	The bitrate announced for this scope.
Up stream port link layer address n	2n	2n[47:0]	48	The link layer to the interface which finishes off the scope of this announce fragment
Down stream port link layer address n	2n+1	(2n+1)[47:0]	48	The link layer address to the interface which starts of the scope of this announce fragment

Table 18: Resource Announce fields

Resource Announce (Static ownership)

This is used if the interfaces ownership distribution is static

0	3	N	broadcast mac addr
0			Announce Source mac
0		Own start	0 Own range
0	br hw 1		Up stream port link layer address 1
	br lw 1		Down stream port link layer address 1
0	br hw n		Up stream port link layer address n
	br lw n		Down stream port link layer address n
0	br hw N		Up stream port link layer address N
	br lw N		Down stream port link layer address N

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56][47:0]	56	Header, Cmd=3
N	0	[55:48]	8	This is the number of fragments in this announce message
Announce Source mac	1	[47:0]	48	This is the mac of the announcing interface
Own Start	2	[47:32]	24	The start of the ownership of the announcer
Own Range	2	[23:0]	24	The range of the ownership of the announcer
br n	2n+1, 2n+2	(2n+1)[55:48] (2n+2)[63:49]	24	The bitrate announced for this scope.
Up stream port link layer address n	2n+1	(2n+1)[47:0]	48	The link layer to the interface which finishes off the scope of this announce fragment
Down stream port link layer address n	2n+2	(2n+2)[47:0]	48	The link layer address to the interface which starts of the scope of this announce fragment

Table 19: Resource Announce fields

Access token passing

These messages are used to request or change access rights for tokens.

5 Resource Request

This message is used to request a set of resource from one interface to another.

0	4	0	Destination mac Address
0		br hw	Source port link layer address
br lw			downstream port link layer address
0			Session identifier

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56][47:0]	56	Header, Cmd=4
Source port link layer address	1	[47:0]	48	The address of the interface which requests to lend resources
requested br	1,2	1[55:48] 2[63:49]	24	The amount of bitrate requested
downstream port link layer address	2	[47:0]	48	The scope for the request, note that the start scope is implicit from the source interface address, since no node will need to borrow resources which does not start at its own scope
B	3	[31:31]	1	If set to one, we indicate that this node is capable of accepting Resource transfers that are multi part.
RtReq session identifier	3	[30:0]	31	An identifier which is to be sent back in the reply to the borrowing request, it is used by the receiver to distinguish between possibly many outstanding borrowing sessions, it needs only be unique for each allocation domain

Table 20: Resource request fields

Resource Transfer

- This is used to transfer the access right from one interface to another. The source of the sent resource is not sent since it is not needed until we return the resources, but by then the owner may well have changed for all or part of the slots. A session identifier is used to map together the request with the reply in case more than one channel is being requested for borrowing at the same time. The session id is also what holds the information on what scope the resource transfer is valid for. In the example below we pad the last slot fragment to the right with 32 bits of zeros.

0	5	N	Destination mac Address	
amt 1	start slot 1		amt 2	start slot 2
amt n-1	start slot n-1		amt n	start slot n
amt N-1	start slot N-1		amt N	start slot N
0			Session identifier	

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56][47:0]	56	Header, Cmd=5
Number of fragments	0	[55:48]	8	The number of fragments in this resource transfer
amt n	n/2+1	(n/2+1) [32*(n%2)+31: 32*(n%2)+24]	8	The slot amount for this fragment
start slot n	n/2+1	(n/2+1) [32*(n%2)+23: 32*(n%2)]	24	The start slot for this fragment
B	N/2+1	[31:31]	1	More to come Bit, this indicates that there will be more slots in a coming message.
RtReq session identifier	N/2+1	[30:0]	31	The session identifier provided by the borrower

Table 21: Resource Transfer fields.

Resource Transfer Return

- 5 This is used to return resources. Again, the owner does not care about nor does it keep records of who originally borrowed the resources. Note the scope fields needed in this message. Below the number of slot fragments is odd we stuff the last (rightmost) 32 bits with zeros.

0	6	N	Destination mac Address
0			Up stream port link layer address
0			Down stream port link layer address
amt 1	start slot 1	amt 2	start slot 2
amt n	start slot n	amt n-1	start slot n-1
amt N-1	start slot N-1	amt N	start slot N

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56][47:0]	56	Header, Cmd=6
Number of fragments	0	[55:48]	8	Number of fragments in Resource Transfer Return
Up stream port link layer address	1	[47:0]	48	The start scope for the return of this resource (applies to all slot fragments)
Down stream port link layer address	2	[47:0]	48	The end scope for the return of this resource (applies to all slot fragments)
amt	$(n-1)/2+3$	$((n-1)/2+3)$ [32*(n%2)+31: 32*(n%2)+24]	8	Slot amount of current fragment
start slot	$(n-1)/2+3$	$((n-1)/2+3)$ [32*(n%2)+23: 32*(n%2)]	24	Slot start of current fragment

Table 22: Resource Transfer Return fields

Probe

- The message responsible for sending out a question for one or more of the owned resources to check if they are still there. The ownership of the resources can be either statically defined in each node for each interface, or it can be negotiated via the dynamic ownership machine. In the dynamic case, it is not OK to send out probes at all times. Please refer to section 12.6.2.5.1 for the exact conditions in which we are allowed to send out probes. The receiver of this message must handle the case when unknown mac addresses come in to the system.

0	7	0	Destination mac Address
0			Source port link layer address
0			Upstream port link layer address
0			Downstream port link layer address
0			Probe sess id
		amount	start slot

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56][47:0]	56	Header, Cmd=7
Source port link layer address	1	[47:0]	48	Interface address of the interface probing for this set of resources
Upstream port link layer address	2	[47:0]	48	The start of the scope for which the probe is valid
Downstream port link layer address	3	[47:0]	48	The end of the scope for which the probe is valid
Probe sess id	4	[47:32]	16	An identifier mostly used for robustness against protocol faults such as multiple outstanding probes for same resource, needs to be unique per allocation domain
amount	4	[31:24]	8	Slot amount of current probed fragment
start slot	4	[23:0]	24	Start of current probed fragment

Table 23: Probe fields

Probe Reply

- 5 This is an answer to a request for resources. Note especially that only the downstream scope of the resources are used in the fields below, the upstream scope is implicit from the source of

this interface (remember that nodes are very unlikely to access borrowed resources that are upstream to them). Note that it is only the resources that are in used (borrowed) that are listed in this reply. During ownership transitions it is possible to get probes for slots that an interface owns itself. It is important to note that if we get a probe that contains queries for slots that are currently not in our ownership range, we must not answer that probe at all.

- 5 For the understanding of this message it is important to note that it is two-dimensional. That is, first we have a number of fragments, one for each active scope. Then we have a set of slot fragments valid for each scope. If in the figure 20 the number of slot fragments, M, is odd for one or any of the slot fragments, 32 bits of zero must be stuffed in the right hand 32 bits of that slot. It is not allowed to send out probe replies at all times, please refer to the earlier description for the exact conditions on when to reply to probes.

0	8	N	Destination mac Address	
probe sess id			Source port link layer address	
0	frags 1	downstream port link layer address 1		
Amt 1	Start 1		Amt 2	Start 2
Amt m	Start m		Amt M	Start M
0	frags 2	downstream port link layer address 2		
Amt 1	Start 1		Amt 2	Start 2
Amt m	Start m		Amt M	Start M
0	frags n	downstream port link layer address n		
Amt 1	Start 1		Amt 2	Start 2
Amt m	Start m		Amt M	Start M
0	frags N	downstream port link layer address N		
Amt 1	Start 1		Amt 2	Start 2
Amt m	Start m		Amt M	Start M

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56][47:0]	56	Header, Cmd=8
N	0	[55:48]	8	Number of different scopes in probe reply
probe sess id	1	[63:48]	16	The id of this probe session, the receiver of this message should check so that this is not old
Source port link layer address	1	[47:0]	48	The source of the interface answering to this probe (remember that probe is broadcast so we need to distinguish the answers, which are unicast)
frags n	f(n,1)	f(n,1)[55:48]	8	Number of fragments for a certain scope in a probe reply
downstream port link layer address n	f(n,1)	f(n,1)[47:0]	48	Down stream address of this probe (upstream is implicit, since no node will use borrowed slots upstream of its own location).
Amt n,m		$f(n,m)[(m\%2)*32+31:(m\%2)*32+24]$	8	The amount of slots in the fragment
Start n,m		$f(n,m)[(m\%2)*32+23:(m\%2)*32]$	24	The start of this slot fragment

Table 24: Probe reply fields.

The function $f(n)$ in table 24 is defined by the formula below. It gives us the slot offset as a function of n and m .

5 This is the function $f(n,m)$ which defines the slot index for various n 's and m 's. Note the variable M_n , which is defined by the number of slot fragments for each scope. We need to use

the sum formula since all the M_n 's may have different values.

Also, if $n=1$, the sum evaluates to zero.

$$f(n,m) = 2 + \frac{m+1}{2} + \sum_{i=1}^{n-1} \left(\frac{M_{m+1}}{2} + 1 \right)$$

Ownership passing

- 5 These are the messages for negotiating ownerships. The starvation bit is a sort of "vertical" ownership information. When receiving a set starvation bit in a sync message, the interface receiving the valid sync message must drop all free resources it may have on that physical link. The interface put in starvation mode must not probe nor respond to probes when it is put in starvation mode via the operators interface.
- 10

Sync

- This message is originated from the master of an allocation domain. The master waits a certain time for the message to reach back again. If it does not reach back, the master re-transmits.
- 15

0	9	N	Destination mac Address	
0			interface link layer address 1	
0			interface link layer address n	
0			interface link layer address N	
0	req own 1		0	req own 2
0	req own n-1		0	req own n
0	req own N-1		0	req own N

- If the number of interfaces in the allocation domain is odd we pad the lower-leftmost 24 bits with zeros.
- 20

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56] [47:0]	56	Header, Cmd=9
Number of interfaces in this allocation domain	0	[55:48]	8	This is the number of interfaces listed in this message, it is also equal to the number of interfaces this interface considered to have when the message was transmitted
S	n	48	1	This indicates that the specific interface does not accept channels in, out or through it
interface link layer address n	n	n[47:0]	48	An interface mac address in the list
req own n	(n-1)/2+3	(n-1)/2+3 [32*(n%2)+23 : 32*(n%2)]	24	The requested ownership limit that this interface wanted last time the master got that information

Table 25: Sync fields

0
1
2
3
4
5
6
7
8
9
A
B
C
D
E
F
G
H
I
J
K
L
M
N
O
P
Q
R
S
T
U
V
W
X
Y
Z
[
,
-
.
/

Gather

The gather is sent by any interface (including the master that then talks to itself) whenever a need to change the policing parameter arises.

5

0	10	0	Destination mac Address
0			interface link layer address
0			req. own

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56] [47:0]	56	Header, Cmd=10
S	1	48	1	Starvation bit, this informs the master that we do not want any more new channels to/from or over this interface.
interface link layer address	1	[47:0]	48	The link layer address of the interface doing the request for a certain ownership range
req. own	2	[23:0]	24	The ownership range requested by this interface

Table 26: Gather fields

Kill

- 10 The owner of a resource issues this message after doing a probe, which detects double booking.

0	11	0	Destination mac Address
0			Kill slot

Fields	Slot #	BitVec	Size	Description
DRMP generic Header	0	[63:56] [47:0]	56	Header, Cmd=11
Kill slot	1	[23:0]	24	The slot requested to be killed due to double booking

Table 27: Kill fields

Questions on system limits.Worst case message sizes

5 A calculation on this needs to be done on a per implementation basis. The parameters defining the message sizes are the following:

- Fragmentation of resources
- Number of interfaces in an allocation domain.

10 It is out of the scope of this document to specify a maximum size of ctrl messages; implementations should look this up and adjust their messages accordingly.

CPU intensity and scalability

15 Unless one has unlimited CPU power it is advisable to implement some kind of CPU meter. Overload will occur in two typical situations:

- An interface, a node or a whole system is coming up and wants to quickly acquire all its resources. This will generate a lot of probe and probe reply messages.
- 20 - A node is having a lot of channel setup and/or tear down going through it, remember that although the network switching is solely done through hardware signaling still can be very costly and CPU overload may occur, especially in nodes centrally placed in the network.

25 In actual implementations it has been found that in many cases the best thing to do is to make each successful probe trigger another one. Then when all resources have been counted once, we go to a more slow recovery probe since message losses during normal operation are relatively rare in their occurrence.

Message reordering

- The DCC protocol is assumed to maintain strict FIFO on the message in all normal cases of operation. Message losses are acceptable at all times, but should be kept to a minimum for performance reasons. Reordering is assumed to be very uncommon, but the kill message will handle the cases that still can occur.

Scalability of the protocol itself

- The probing process can be quite costly for the nodes involved. However the cost is linear or near linear provided that there is an inexpensive way to filter incoming messages not directed to the node itself.

Consider the process of bootstrapping a system with N nodes and M slots.

- Assume the largest possible message size of the system corresponds to k probe replies in the same message
Assume the ownerships (and thus also the obligation to probe for resources) is distributed even among the nodes.
- This discussion is valid for $N \geq 2$.
 - M and k are independent of N .
 - Each node sends out $\frac{M}{kN}$ probes.
 - A probe consists of $N-1$ queries and $N-1$ replies.
 - This gives us a total of $\frac{2M(N-1)}{kN}$ messages sent in the system.
 - This can be rewritten as $\frac{2M}{k}(1 - \frac{1}{N})$ and $\frac{1}{N}$ diminishes towards zero as N increases.
 - This implies that the cost per node is equal to or smaller than $\frac{2M}{k}$ i.e. nearly constant as N grows.

- Since the number of nodes increase when one employs a larger bypass chain the processing power also increases in the total system, but so does the number of processors sharing the work. This shows that the CPU cost of the system is constant if each node only gets the messages intended for it.

If filtering of messages not intended for a specific recipient is associated with a non-neglectable cost, an additional cost proportional to N is introduced. This is always the case for bitrate cost.

5 Probe optimizations

Here is a set of hints and descriptions on how to make probing more efficient.

Turbo probe

- 10 An optimization concept aiming at having a faster probe rate when we startup or have detected inconsistencies and a slower probe rate during normal operation. This is done to minimize the bootup time for a node. The node has just started. It always starts with no resources but with an ownership range assigned. Thus it also has a responsibility to probe its resource tokens at
- 15 a regular interval. At bootup, the probe is intensified. The turbo strategy is to make at least one successful probing for each resource unit and when that is done go to probing at lower rate (and thus spend less CPU). The turbo is also re triggered whenever the topology of the specific allocation domain changes
- 20 (due to fiber or node failures). When this initial probing has been completed we settle for a relatively slow mechanism for the detection of lost resources, since the probability of message losses in a DTM system is relatively low.

Re-triggered-on-reply probe

- 25 The turbo probe is the concept of sending probes at a faster rate when needed. The re-triggered-on-reply probe is a flow rate mechanism for that. Generally it is hard to make a smooth send-out of probes when using a timer. The idea of re-triggered-on-reply is to let the last probe reply message for a given resource
- 30 or resource-set trigger the next transmission of a probe. Sending probes periodically on a timer seems like a good idea, but in a real implementation case it often results in packet bursts since the granularity of most operating systems clocks is quite coarse. Scheduling only a few hundred timers a second
- 35 and therefore also a few hundred context switches would be

noticeable on the CPU meter even when not probing at all. However with a turnaround from the querying interface to the answering interfaces of around 1 ms we could easily achieve 1000 probes per second without buffer build-up problems. It has
 5 been argued that this will not work for longer networks, but it is unlikely that anyone will build anything else than point-to-point in a very long link (such as a transatlantic or nation-wide one). The advantage of the ptp is that since the allocation domains are local (one on each side), no messages need to be
 10 exchanged with the other side in order to retrieve the resources.

Naturally, all of the process steps, as well as any sub-sequence of steps, described above may be carried out by means of a computer program being directly loadable into the internal memory of a computer, which includes appropriate software for
 15 performing the necessary steps when the program is run on a computer. The computer program can likewise be recorded onto arbitrary kind of computer readable medium.

Generally, a range of the ownership to the resources for a particular interface are distributed according to the expression:
 20 $\lfloor (V_{\max} \times P) / (\Sigma_{\text{req}}) \rfloor$;

where

V_{\max} denotes a maximal number of resources in the system,

P denotes a number of resources requested by the particular interface, and

25 Σ_{req} denotes a sum of the number of resources requested by all interfaces.

However, in order to elucidate the practical consequences of the proposed method of allocating resources, three different calculation examples are shown below.

30 In a first example the system is presumed to have 2000 slots in total, i.e. $V_{\max} = 2000$. Four interfaces A, B, C and D all request 2000 slots each, i.e. $P = 2000$. An equal range of

$\lfloor (2000 \times 2000)/8000 \rfloor = 500$ is thus allocated to all the interfaces A, B, C and D.

In case (i) such amounts of resources are requested by the interfaces in the system that the quotient $P/(\Sigma_{req})$ results in a fractional number, and (ii) a total number of resources has been requested by the interfaces A, B, C and D, which is larger than the number of resources in the system, i.e. $1/(\Sigma_{req}) > V_{max}$, the ownership to any surplus resources is allocated to the master interface according to the expression: $V_{max} - \Sigma_{alloc}$; where Σ_{alloc} denotes a sum of ranges already allocated to the interfaces, i.e. the master interface as well as one or more slave interfaces.

If however, less than the total number of resources in the system have been requested by the interfaces A, B, C and D, or exactly the total number of resources has been requested (i.e. $\Sigma_{req} \leq V_{max}$) no additional resources are allocated to the master interface.

In a second example the system is again presumed to have 2000 slots in total, i.e. $V_{max} = 2000$. A first interface A requests 1700 slots, i.e. $P_A = 1700$; a second interface B requests 2000 slots, i.e. $P_B = 2000$; a third interface C requests 500 slots, i.e. $P_C = 500$; and a fourth interface D requests 700 slots, i.e. $P_D = 700$. The third interface C is presumed to be the master interface. Now, the ownership ranges of the resources are distributed according to the following.

25 A: $\lfloor (2000 \times 1700)/4900 \rfloor = 693$,
 B: $\lfloor (2000 \times 2000)/4900 \rfloor = 816$,
 C: $\lfloor (2000 \times 500)/4900 \rfloor = 204$, and
 D: $\lfloor (2000 \times 700)/4900 \rfloor = 285$.

$\Sigma_{alloc} = 693 + 816 + 204 + 285 = 1998$.
 30 i.e. $V_{max} - \Sigma_{alloc} = 2000 - 1998 = 2$. This means that the master interface C is allocated the ownership of an additional 2 slots,

thus in total $204 + 2 = 206$ slots.

Finally, we show an example where less than the total number of slots are requested by the interfaces in the system. Again 2000 slots are presumed to be available, i.e. $V_{\max} = 2000$. A first interface A requests 200 slots, i.e. $P_A = 200$; a second interface B requests 2000 slots, i.e. $P_B = 300$; a third interface C (master) requests 100 slots, i.e. $P_C = 100$; and a fourth interface D requests 50 slots, i.e. $P_D = 50$. Since here, $\Sigma_{\text{alloc}} < V_{\max}$ each interface is allocated the ownership to exactly as many resources as the respective interface has requested. Thus, the ownership ranges of the resources are distributed according to the following.

A: 200, B: 300, C: 100 and D: 50.

The remaining 1350 are unreserved and held available for use by any future interfaces in the system.

The actual resources owned by each interface is calculated by the respective interface on basis of the particular interface's topological position relative to the range of the ownership to the resources allocated to the interface. In practice, this means that a first interface in a sequence, for instance the first interface A in the example above, is allocated the ownership to the first resources in accordance with its range, i.e. $1 \rightarrow 200$. Correspondingly, a following interface, for instance the second interface B, is allocated the ownership to the next range of resources, i.e. $201 \rightarrow (201 + 300) = 201 \rightarrow 501$, a yet following interface, for instance the third interface C, is allocated the ownership to a following range of resources, i.e. $502 \rightarrow (502 + 100) = 502 \rightarrow 602$ and a last interface, for instance the fourth interface D, is allocated the ownership to the range of resources, i.e. $603 \rightarrow (603 + 50) = 603 \rightarrow 653$.

In general terms, this may be expressed as:

$$\Sigma tp_{\text{lower}} + 1 \rightarrow \Sigma tp_{\text{lower}} + R_i$$

where Σtp_{lower} denotes a sum of ranges for all interfaces having a lower topological position number than the particular interface, and R_i denotes the range for the particular interface.

5 The term "comprises/comprising" when used in this specification is taken to specify the presence of stated features, integers, steps or components. However, the term does not preclude the presence or addition of one or more additional features, integers, steps or components or groups thereof.

10 The invention is not restricted to the described embodiments in the figures, but may be varied freely within the scope of the claims.

9
1
2
3
4
5
6
7
8

Claims

1. A method of allocating resources in a synchronous time division multiplex communications system having a master interface communicating with one or more slave interfaces, and
5 in which the resources between the interfaces are represented by time slots in a repeating frame structure, characterized by the steps of:
- 10 sending a link status message from an interface whenever the interface registers a change in the topology of the system,
 - 10 sending a gather message from an interface whenever the interface requests a revision of a current ownership distribution of resources,
 - 15 sending a sync message from the master interface as an indication of a current distribution of ownership with respect to the resources between the interfaces in the system, and
 - 15 for each interface generating a distribution of the ownership to the resources on basis of the interface's topological position and a latest received sync message.
2. A method according to claim 1, characterized by sending
20 the link status message to all interfaces in the system.
3. A method according to any one of the preceding claims, characterized by sending the gather message to the master interface.
4. A method according to any one of the preceding claims,
25 characterized by sending the sync message to all interfaces in the system, the sync message including information pertaining to a number of resources requested by the respective interfaces in the system, and the sync message being updated until all interfaces refrain from initiating any further gather messages.

5. A method according to any one of the preceding claims, characterized by generating a distribution range of the ownership to the resources for a particular interface according to: $\lfloor (V_{\max} \times P) / (\Sigma_{\text{req}}) \rfloor$;

5 where

V_{\max} denotes a maximal number of resources in the system,

P denotes a number of resources requested by the particular interface, and

10 Σ_{req} denotes a sum of the number of resources requested by all interfaces.

6. A method according to claim 5, characterized by allocating an additional range of the ownership to the resources to the master interface according to:

15
$$V_{\max} - \Sigma_{\text{alloc}}; \quad \text{if } \Sigma_{\text{req}} > V_{\max}, \text{ and}$$

$$0; \quad \text{if } \Sigma_{\text{req}} \leq V_{\max}$$

where Σ_{alloc} denotes a sum of ranges allocated to the master interface and the at least one slave interface.

20 7. A method according to any one of the claims 5 or 6, characterized by allocating ownership to the resources for a particular interface with respect to the interface's topological position relative to the range of the ownership to the resources allocated to the interface.

25 8. A method according to claim 7, characterized by allocating ownership to the resources for a particular interface having a range: $\Sigma_{\text{tp}_{\text{lower}}} + 1 \rightarrow \Sigma_{\text{tp}_{\text{lower}}} + R_i$

where

$\Sigma_{\text{tp}_{\text{lower}}}$ denotes a sum of ranges for all interfaces having a lower topological position number than the particular

interface, and

R_i denotes the range for the particular interface.

9. A computer program directly loadable into the internal memory of a computer, comprising software for performing the method according to of any one of the claims 1 – 8 when said program is run on the computer.

10. A computer readable medium, having a program recorded thereon, where the program is to make a computer perform the method according to of any one of the claims 1 – 8.

10 11. A communications system having transmission resources in the form of time slots in a repeating frame structure, in which the time slots are dynamically allocable, the system comprising at least two interfaces of which one is a master interface and at least one is a slave interface, characterized in that it comprises at least one node, which in turn includes one or more of the interfaces, the node being adapted to effecting the method of according to any one of the claims 1 – 8.

0
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

Abstract

The present invention relates to a dynamic allocation of resources in a synchronous time division multiplex communications system. The system is presumed to include a master interface, which communicates with one or more slave interfaces. Furthermore, the resources are represented by time slots in a repeating frame structure, such as in a DTM system. The invention involves sending a link status message from an interface whenever the interface registers a change in the topology of the system. A gather message is sent from a particular interface to the master interface whenever the interface in question requests a revision of a current ownership distribution of resources. The master interface sends a sync message as an indication of a current distribution of ownership with respect to the resources between the interfaces in the system. Each interface generates a distribution of the ownership to the resources on basis of the interface's topological position and a latest received sync message.

9
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192

1/5

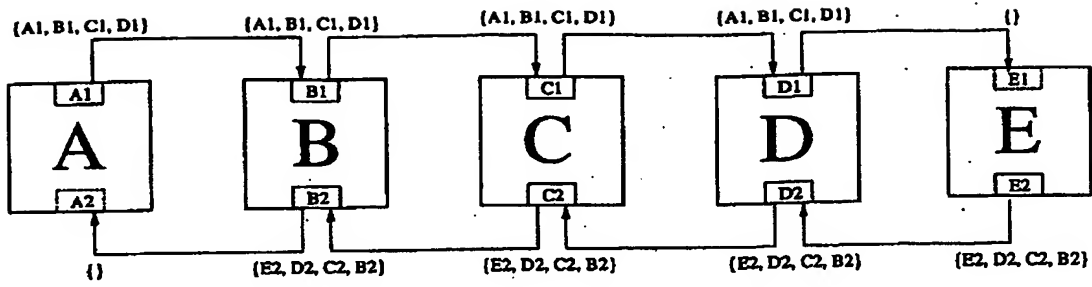


Fig. 1a

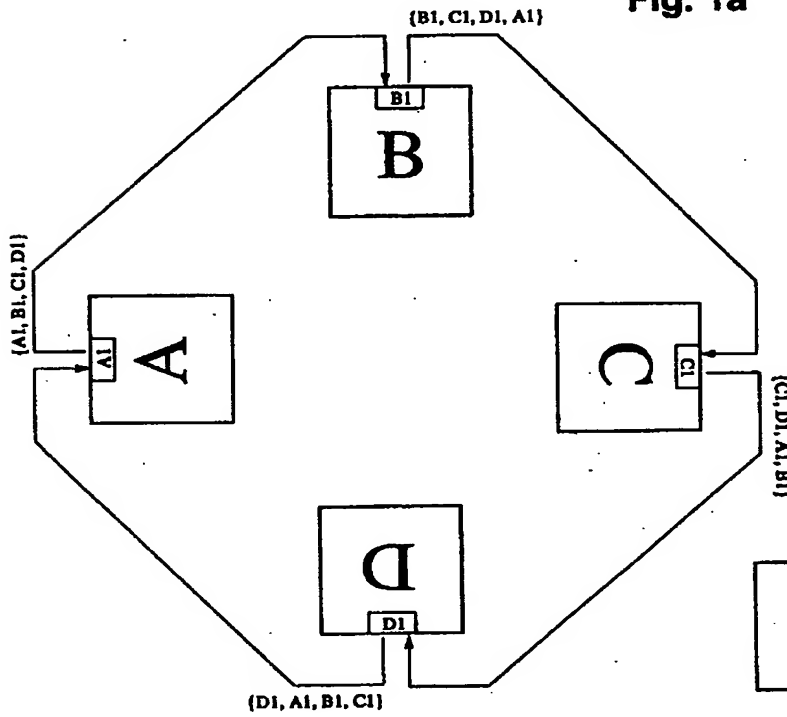


Fig. 1b

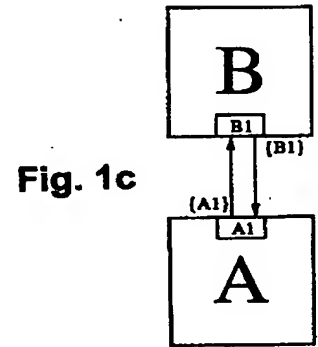


Fig. 1c

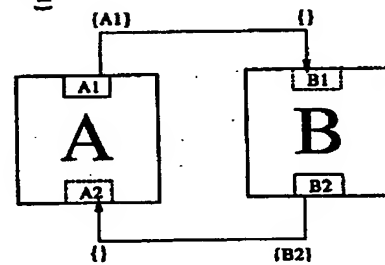
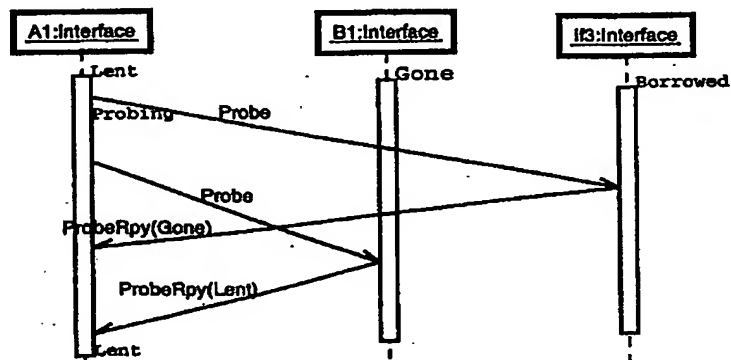


Fig. 1d



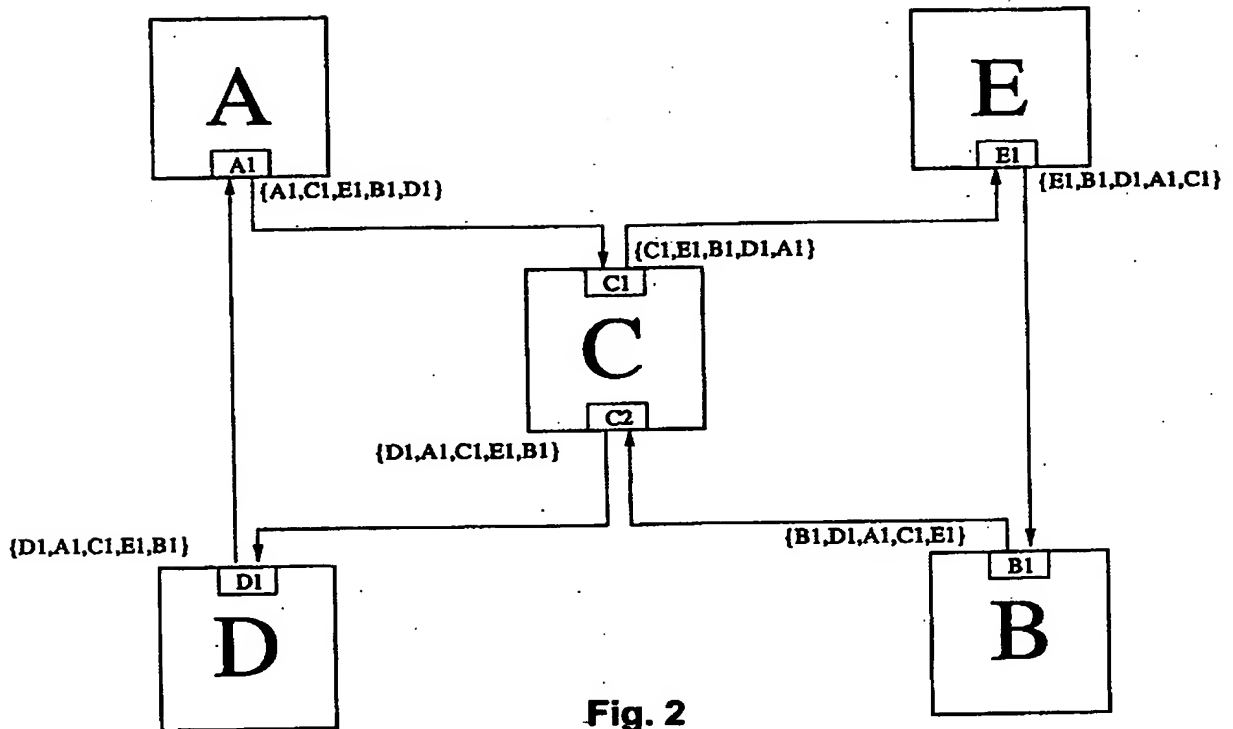


Fig. 2

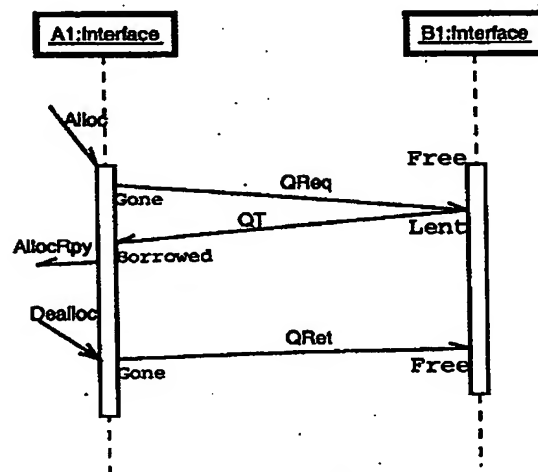
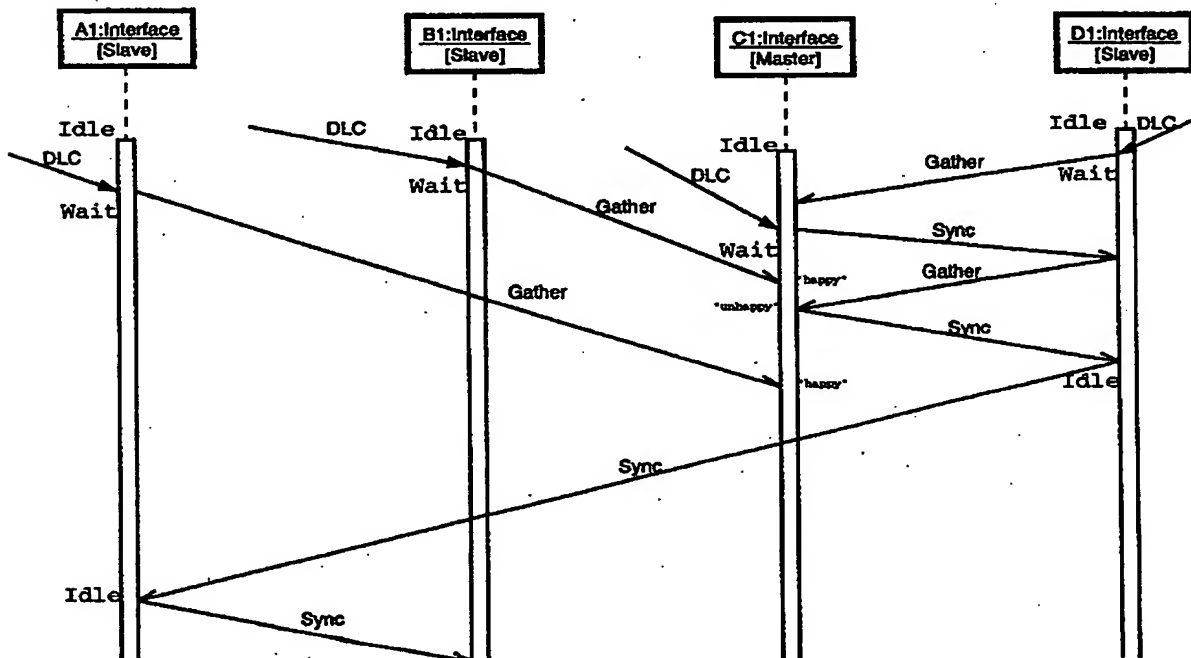
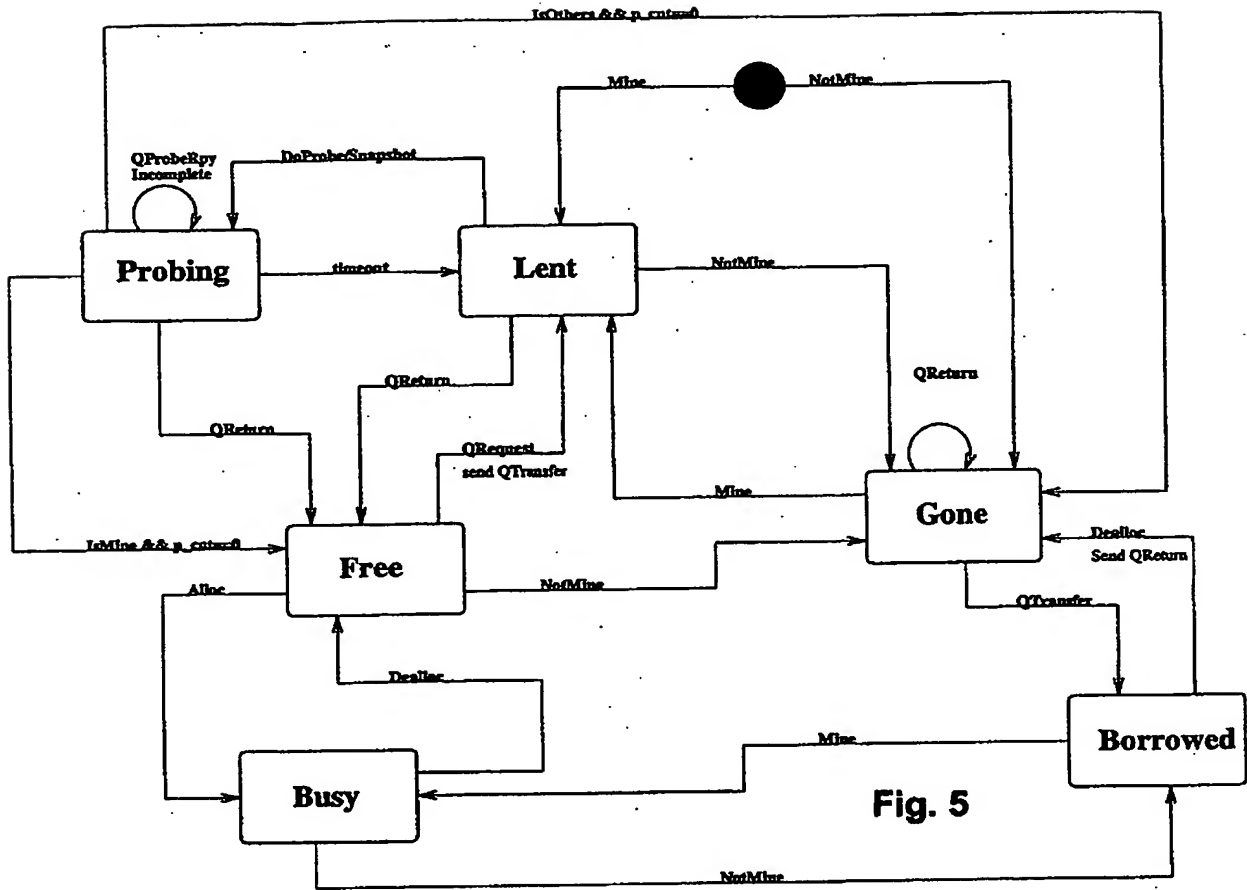


Fig. 4

3/5



4/5

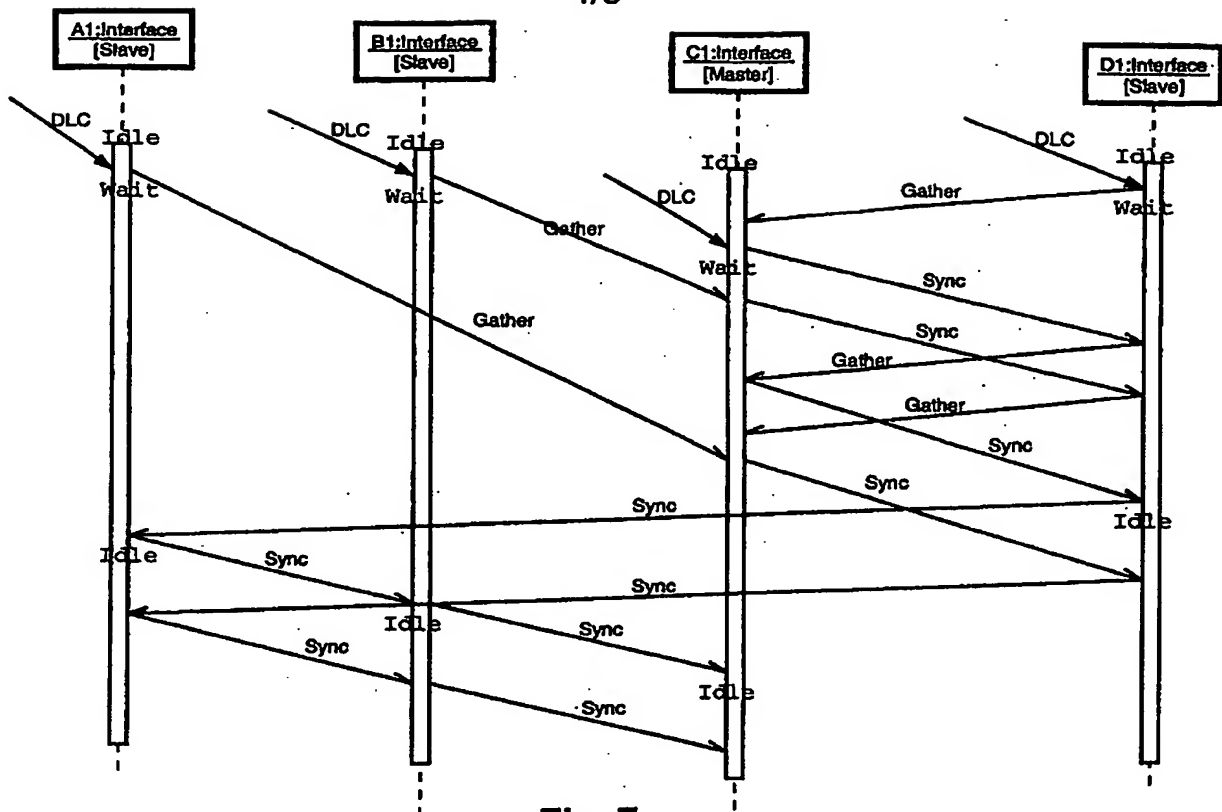
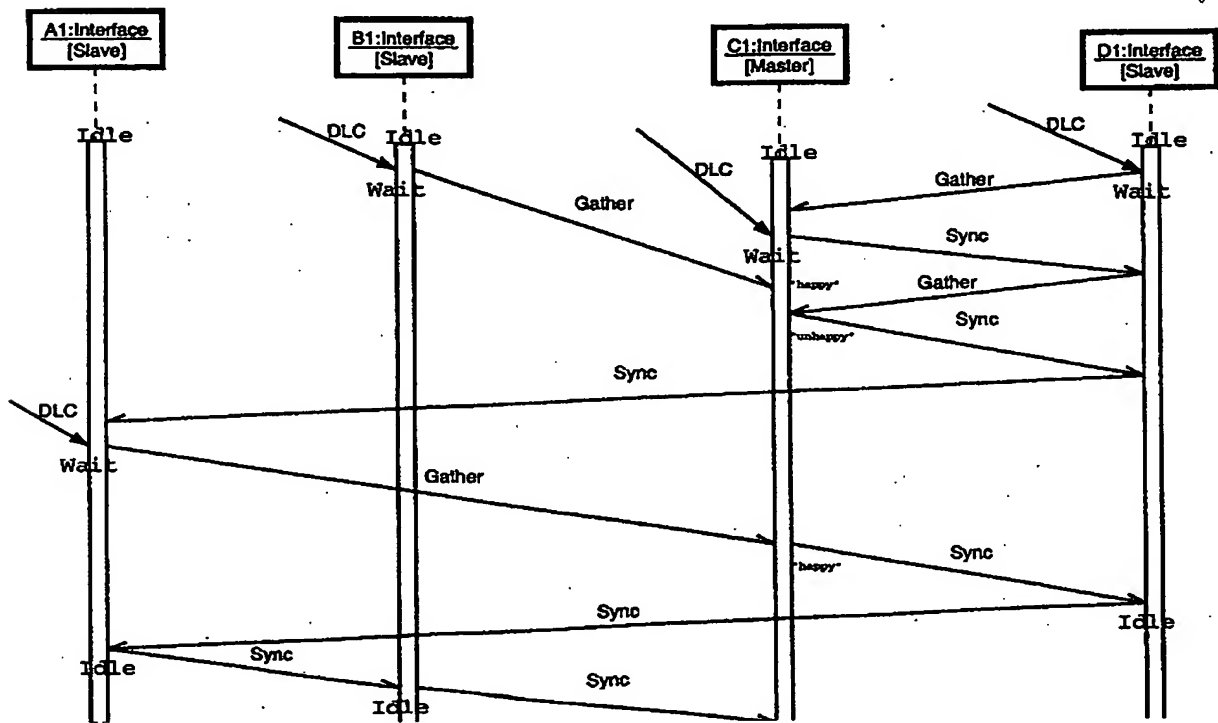


Fig. 7



5/5

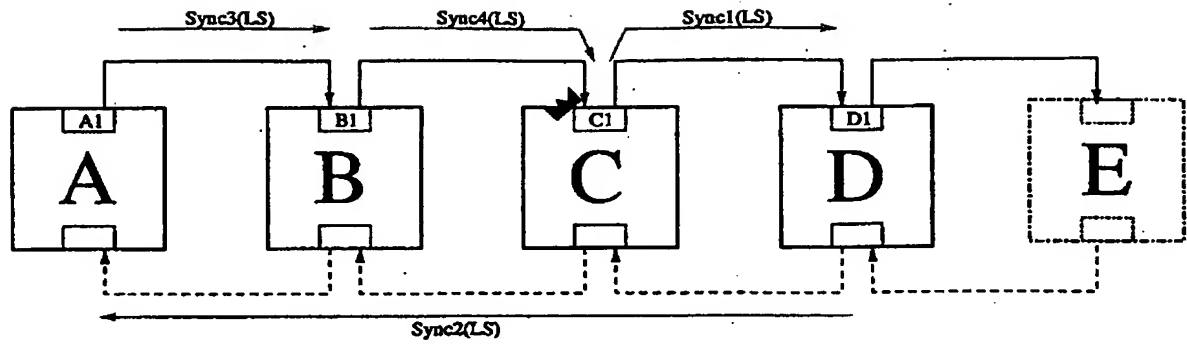


Fig. 9a

AD = { A1, B1, C1, D1 }

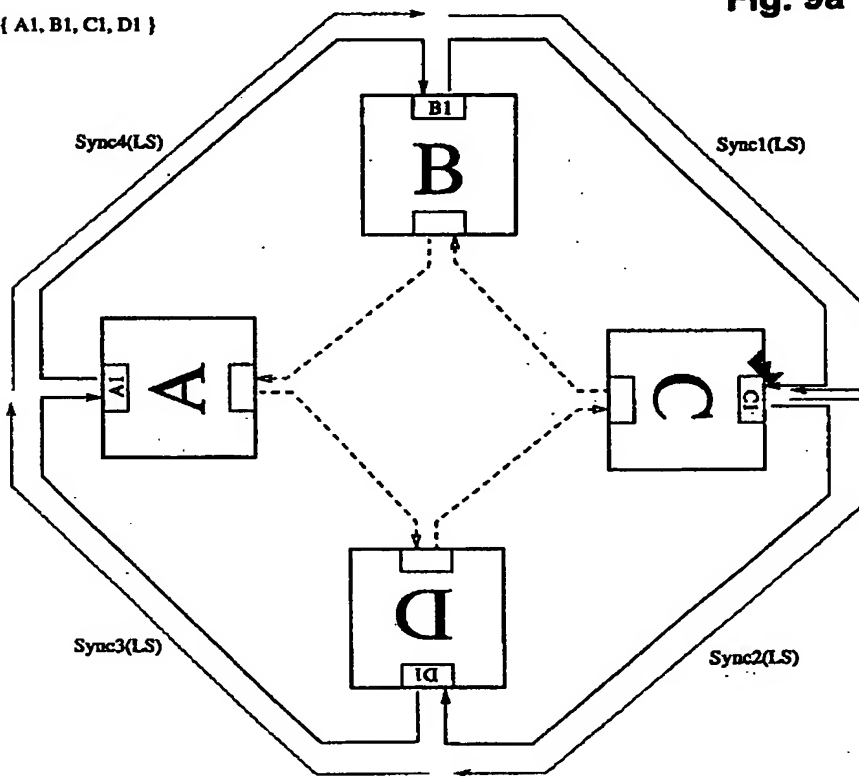


Fig. 9b